

AD-A186 712 VARIANCE FUNCTION ESTIMATION REVISION(U) NORTH CAROLINA 1/1

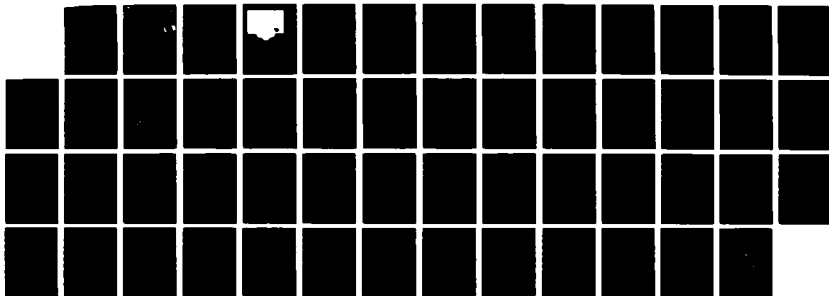
UNIV AT CHAPEL HILL INST OF STATISTICS

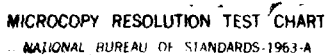
M DAVIDIAN ET AL. MAR 87 MIMED SER-1700-REV

UNCLASSIFIED AFOSR-TR-87-1102 F49620-85-C-0144

F/G 12/3

NL





MICROCOPY RESOLUTION TEST CHART  
NATIONAL BUREAU OF STANDARDS-1963-A

OTIC FILE COPY

UNCLASSIFIED

SECURITY CLASSIFICATION

REPORT DOCUMENTATION PAGE

AD-A186 712

1b. RESTRICTIVE MARKINGS

3. DISTRIBUTION/AVAILABILITY OF REPORT  
Approved for public release; distribution unlimited.

4. PERFORMING ORGANIZATION REPORT NUMBER(S)

Mimeo Series #1700

5. MONITORING ORGANIZATION REPORT NUMBER(S)  
AFOSR-TR-87-1102

6a. NAME OF PERFORMING ORGANIZATION

University of North Carolina

6b. OFFICE SYMBOL  
(If applicable)

7a. NAME OF MONITORING ORGANIZATION

Air Force Office of Scientific Research

6c. ADDRESS (City, State and ZIP Code)

Stat. Dept, UNC-CH  
Chapel Hill, NC 27514

7b. ADDRESS (City, State and ZIP Code)

Worm Creek 8C

8a. NAME OF FUNDING/SPONSORING ORGANIZATION

AFOSR

8b. OFFICE SYMBOL  
(If applicable)

nm

9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER

149620-85-C-0144

8c. ADDRESS (City, State and ZIP Code)

Bolling Air Force Base  
Washington, DC 20332

10. SOURCE OF FUNDING NOS.

PROGRAM  
ELEMENT NO.

PROJECT  
NO.

TASK  
NO.

WORK UNIT  
NO.

611021 2304 AS

11. TITLE (Include Security Classification)

"Variance Function Estimation (revised)"

12. PERSONAL AUTHOR(S)

Davidian, Marie and Carroll, R.J.

13a. TYPE OF REPORT

technical journal

13b. TIME COVERED

FROM 8/86 TO 8/87

14. DATE OF REPORT (Yr., Mo., Day)

March 1987

15. PAGE COUNT

42

16. SUPPLEMENTARY NOTATION

17. COSATI CODES

FIELD	GROUP	SUB. GR.

18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number)

Heteroscedastic regression models, standard asymptotic theory, variance function estimation

19. ABSTRACT (Continue on reverse if necessary and identify by block number)

Heteroscedastic regression models are used in fields including economics, engineering, and the biological and physical sciences. Often, the heteroscedasticity is modeled as a function of the covariates or the regression and other structural parameters. Standard asymptotic theory implies that how one estimates the variance function, in particular the structural parameters, has no effect on the first order properties of the regression parameter estimates; however, there is evidence both in practice and higher order theory to suggest that how one estimates the variance function does matter. Further, in some settings estimation of the variance function is of independent interest or plays an important role in estimation of other quantities. In this paper, we study variance function estimation in a unified way, focusing on common methods proposed in the statistical and other literature, in order to make both general observations and compare different estimation schemes. We show there are significant differences in both efficiency and robustness for many common methods. (look on back)

20. DISTRIBUTION/AVAILABILITY OF ABSTRACT

UNCLASSIFIED/UNLIMITED ☒ SAME AS RPT. ☐ OTIC USERS ☐

21. ABSTRACT SECURITY CLASSIFICATION

UNCLASSIFIED

22a. NAME OF RESPONSIBLE INDIVIDUAL

Lisa Brooks

22b. TELEPHONE NUMBER (Include Area Code)

919-962-2707

22c. OFFICE SYMBOL

nm

DTIC  
ELECTE  
OCT 07 1987

*s. de... col*

We develop a general theory for variance function estimation, focusing on estimation of the structural parameters and including most methods in common use in our development. The general qualitative conclusions are these. First, most variance function estimation procedures can be looked upon as regressions with "responses" being transformations of absolute residuals from a preliminary fit or sample standard deviations from replicates at a design point. ~~Our conclusion is that~~ the former is typically more efficient, but not uniformly so. Secondly, for variance function estimates based on transformations of absolute residuals, we show that efficiency is a monotone function of the efficiency of the fit from which the residuals are formed, at least for symmetric errors. Our conclusion is that one should iterate so that residuals are based on generalized least squares. Finally, robustness issues are of even more importance here than in estimation of a regression function for the mean. The loss of efficiency of the standard method away from the normal distribution is much more rapid than in the regression problem. ←

As an example of the type of model and estimation methods we consider, for observation - covariance pairs  $(Y_i, x_i)$ , one may model the variance as proportional to a power of the mean response, e.g.,

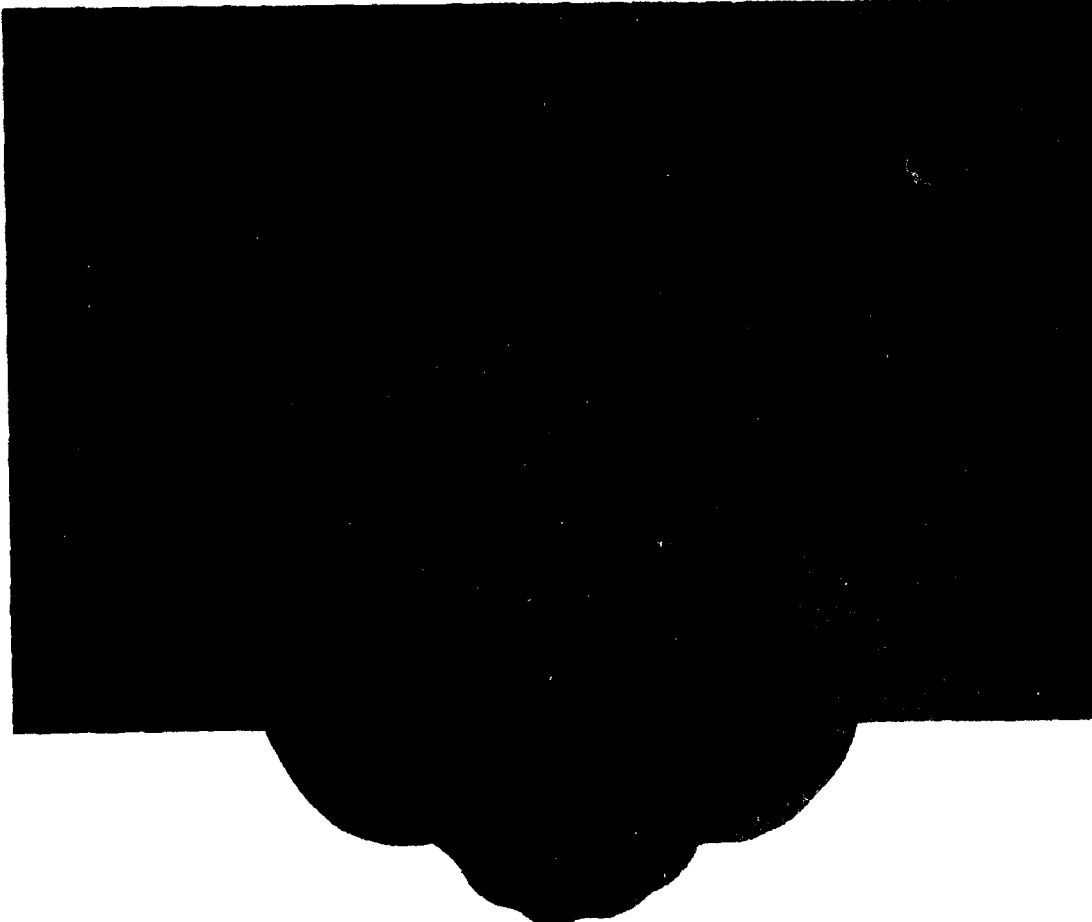
$$E(Y_i) = f(x_i, \beta); \text{var}(Y_i) = \sigma f(x_i, \beta)^\theta, f(x_i, \beta) > 0,$$

where  $f(x_i, \beta)$  is the possibly nonlinear mean function and  $\theta$  is the structural parameter of interest. "Regression methods" for estimation of  $\theta$  and  $\sigma$  based on residuals

$r_i = Y_i - f(x_i, \hat{\beta}_*)$  for some regression fit  $\hat{\beta}_*$  involve minimizing a sum of squares where some function  $T$  of the  $|r_i|$  plays the role of the "responses" and an appropriate function of the variance plays the role of the "regression function." For example, if  $T(x) = x^2$ , the responses would be  $r_i^2$ , and the form of the regression function would be suggested by the approximate fact  $E(r_i^2) \sim \sigma^2 f(x_i, \hat{\beta}_*)^{2\theta}$ . One could weight the sum of squares appropriately by considering the approximate variance of  $r_i^2$ . For the case of replication at each  $x_i$ , some methods suggest replacing the  $r_i$  in the function  $T$  by sample standard deviations at each  $x_i$ . Other functions  $T$ , such as  $T(x) = x$  or  $\log x$  have also been proposed.



Accession For	
NTIS GRA&I	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By _____	
Distribution/	
Availability Codes	
Dist	Avail and/or Special
A-1	



VARIANCE FUNCTION ESTIMATION  
(revised)

by  
Marie Davidian

&

R.J. Carroll

Mimeo Series #1700

March 1987

DEPARTMENT OF STATISTICS  
Chapel Hill, North Carolina

VARIANCE FUNCTION ESTIMATION

M. Davidian and R.J. Carroll

Department of Statistics  
University of North Carolina at Chapel Hill  
321 Phillips Hall 039 A  
Chapel Hill, North Carolina 27514

### ABSTRACT

Heteroscedastic regression models are used in fields including economics, engineering, and the biological and physical sciences. Often, the heteroscedasticity is modeled as a function of the covariates or the regression and other structural parameters. Standard asymptotic theory implies that how one estimates the variance function, in particular the structural parameters, has no effect on the first order properties of the regression parameter estimates; however, there is evidence both in practice and higher order theory to suggest that how one estimates the variance function does matter. Further, in some settings, estimation of the variance function is of independent interest or plays an important role in estimation of other quantities. In this paper, we study variance function estimation in a unified way, focusing on common methods proposed in the statistical and other literature, in order to make both general observations and compare different estimation schemes. We show there are significant differences in both efficiency and robustness for many common methods.

We develop a general theory for variance function estimation, focusing on estimation of the structural parameters and including most methods in common use in our development. The general qualitative conclusions are these. First, most variance function estimation procedures can be looked upon as regressions with "responses" being transformations of absolute residuals from a preliminary fit or sample standard deviations from replicates at a design point. Our conclusion is that the former is typically more efficient, but not uniformly so. Secondly, for variance

function estimates based on transformations of absolute residuals, we show that efficiency is a monotone function of the efficiency of the fit from which the residuals are formed, at least for symmetric errors. Our conclusion is that one should iterate so that residuals are based on generalized least squares. Finally, robustness issues are of even more importance here than in estimation of a regression function for the mean. The loss of efficiency of the standard method away from the normal distribution is much more rapid than in the regression problem.

As an example of the type of model and estimation methods we consider, for observation - covariate pairs  $(Y_i, x_i)$ , one may model the variance as proportional to a power of the mean response, e.g.,

$$E(Y_i) = f(x_i, \beta) \quad ; \quad \text{var}(Y_i) = \sigma f(x_i, \beta)^\theta, \quad f(x_i, \beta) > 0,$$

where  $f(x_i, \beta)$  is the possibly nonlinear mean function and  $\theta$  is the structural parameter of interest. "Regression methods" for estimation of  $\theta$  and  $\sigma$  based on residuals  $r_i = Y_i - f(x_i, \hat{\beta}_*)$  for some regression fit  $\hat{\beta}_*$  involve minimizing a sum of squares where some function  $T$  of the  $|r_i|$  plays the role of the "responses" and an appropriate function of the variance plays the role of the "regression function." For example, if  $T(x) = x^2$ , the responses would be  $r_i^2$ , and the form of the regression function would be suggested by the approximate fact  $E(r_i^2) \approx \sigma^2 f(x_i, \hat{\beta}_*)^{2\theta}$ . One could weight the sum of squares appropriately by considering the approximate variance of  $r_i^2$ . For the case of replication at each  $x_i$ , some methods suggest replacing the  $r_i$  in the function  $T$  by sample standard deviations at each  $x_i$ . Other functions  $T$ , such as  $T(x) = x$  or  $\log x$  have also been proposed.



## 1. INTRODUCTION

Consider a heteroscedastic regression model for observable data  $Y$ :

$$(1.1) \quad EY_1 = \mu_1 = f(x_1, \beta); \quad \text{Var}(Y_1) = \sigma^2 g^2(z_1, \beta, \theta).$$

Here,  $\{x_1\}$  are the design vectors,  $\beta(p \times 1)$  is the regression parameter,  $f$  is the mean response function, and the variance function  $g$  expresses the heteroscedasticity, where  $\{z_1\}$  are known vectors, possibly the  $\{x_1\}$ ,  $\sigma$  is an unknown scale parameter, and  $\theta(r \times 1)$  is an unknown parameter. For example, the variance may be modeled as proportional to a power of the mean:

$$(1.2) \quad g(z_1, \beta, \theta) = f(x_1, \beta)^\theta, \quad f(x_1, \beta) > 0.$$

One might also model the variance as quadratic in the predictors, i.e.,

$$\sigma g(z_1, \beta, \theta) = 1 + \theta_1 x_1 + \theta_2 x_1^2,$$

or by an expanded power of the mean model, i.e.,

$$(1.3) \quad \sigma^2 g^2(z_1, \beta, \theta) = \theta_1 + \theta_2 f(x_1, \beta)^{\theta_3}.$$

Box and Meyer (1986) use

$$g(z_1, \beta, \theta) = \exp(z_1^t \theta).$$

An important feature of (1.1) is that no assumption about the distribution of

the  $\{Y_i\}$  has been made other than that of the form of the first two moments. Models which may be regarded as special cases of (1.1) are used in diverse fields, including radioimmunoassay, econometrics, pharmacokinetic modeling, enzyme kinetics and chemical kinetics among others. The usual emphasis is on estimation of  $\beta$  with estimation of the variances as an adjunct.

The most common method for estimating  $\beta$  is generalized least squares, in which one estimates  $g(z_i, \beta, \theta)$  by using an estimate of  $\theta$  and a preliminary estimate of  $\beta$  and then performs weighted least squares; see, for example, Carroll and Ruppert (1982a) and Box and Hill (1974). This might be iterated, with the preliminary estimate replaced by the current estimate of  $\beta$ , a new estimate of  $\theta$  obtained and the process repeated. Standard asymptotic theory as in Carroll and Ruppert (1982a) or Jobson and Fuller (1980) shows that as long as the preliminary estimators for the parameters of the variance function are consistent, all estimators of  $\beta$  obtained in this way will be asymptotically equivalent to the weighted least squares estimator with known weights.

There is evidence that for finite samples, the better one's estimate of  $\theta$ , the better one's final estimate of  $\beta$ . Williams (1975) states that "both analytic and empirical studies...indicate that...the ordering of efficiency (of estimates of  $\beta$ )...in small samples is in accordance with the ordering by efficiency (of estimates of  $\theta$ ).". Rothenberg (1984) shows via second order calculations that if  $g$  does not depend on  $\beta$ , when the data are normally distributed the covariance matrix of the generalized least squares estimator of  $\beta$  is an increasing function of the covariance matrix of the estimator of  $\theta$ .

Second order asymptotics provide only a weak justification for studying the properties of variance function estimates. Instead, our thesis is that estimation of the structural variance parameter  $\theta$  is of independent interest. In many engineering applications, an important goal is to estimate the error made in predicting a new observation; this can be obtained from the variance

function once a suitable estimate of  $\theta$  is available. In chemical and biological assay problems, issues of prediction and calibration arise. In such problems, the estimator of  $\theta$  plays a central role. As motivation for the study of variance function estimation, in Section 2 we discuss the problems of calibration and prediction in the case of heteroscedasticity. For a formal investigation of how the statistical properties of prediction intervals and calibration constructs such as the minimal detectable concentration will be highly dependent on how one estimates  $\theta$ ; see Carroll, Davidian and Smith (1986). In off-line quality control, the emphasis is not only on the mean response but also on its variability; Box and Meyer (1986) state that "one distinctive feature of Japanese quality control improvement techniques is the use of statistical experimental design to study the effect of a number of factors on variance as well as the mean." The goal is to adjust the levels of a set of experimental factors to bring the mean of the responses to some target value while minimizing standard deviation; the problem involves simultaneous consideration of both mean and variability, where the latter may be a function of the mean, see Box (1986) and Box and Ramirez (1986). These authors advocate methods based on data transformations to account for the heteroscedasticity in separating the factors into those affecting dispersion but not location, those affecting location but not dispersion, and those affecting neither. One similarly might employ effective estimation of variance functions in this application. We briefly discuss the relationship between variance function estimation and one type of data transformation in Section 3.

It should be evident from this brief review that far from being only a nuisance parameter, the structural variance parameter  $\theta$  can be an important part of a statistical analysis. The above discussion suggests the need for a unified investigation of estimation of variance functions, in particular, estimation of the structural parameter  $\theta$ . Previous work in the literature

tends to treat various special cases of (1.1) as different models with their own estimation methods. The intent of this paper is to study parametric variance function estimation in a unified way. Nonparametric variance function estimation has also been studied, see for example Carroll (1982); we will confine our study to the parametric setting.

Parametric variance function estimation may be thought of as a type of regression problem in which we try to understand variance as a function of known or estimable quantities, and in which  $\theta$  plays the part of a "regression" parameter. The major insight which allows for a unified study is that the absolute residuals from the current fit to the mean or the sample standard deviations from replicates are basic building blocks for analysis. At the graphical level, this means that transformations of the absolute residuals and sample standard deviations can be used to gain insight into the structure of the variability and to suggest parametric models. For estimation, a major contribution is to point out that most of the methods proposed in the literature are (possibly weighted) regressions of transformations of the basic building blocks on their expected values. Many exceptions to this are dealt with in this article as well.

Our study yields these major qualitative conclusions. As stated here, they apply strictly only to symmetric error distributions, but they are fairly definitive, and one is unlikely to be too successful ignoring them in practice.

(I). Robustness plays a great role in the efficiency of variance function estimation, probably even greater than in estimation of a mean function. For example, if the variance does not depend on the mean response, the standard method will be normal theory maximum likelihood as in Box & Meyer (1986). A weighted analysis of absolute residuals yields an estimator only 12% less efficient at the normal model which rapidly gains efficiency over maximum likelihood for progressively heavier-tailed distributions. This slope of

improvement is much larger than is typical in regression on means. For a standard contaminated normal model for which the best robust estimators have efficiency 125% with respect to least squares, the absolute residual estimator of the variance function has efficiency 200%.

(II). We obtain implications for fit to the means upon which the residuals are based. It has been our experience that unweighted least squares residuals yield unstable estimates of the variance function when the variances depend on the mean. This is confirmed in our study, in the sense that the asymptotic efficiency of the variance function estimators is an increasing function of the efficiency of the current fit to the means. Thus, we suggest the use of iterative weighted fitting, so that the variance function estimate is based on generalized least squares residuals. As far as we can tell, this part of our paper is one of the first formal justifications for iteration in a generalized least squares context.

(III). It is standard in many applied fields to take  $m$  replicates at each design point, where usually  $m \leq 4$ . Rather than using (transformations of) absolute residuals for estimating variance function parameters, one might use the sample standard deviations. We develop an asymptotic theory from which we obtain the efficiency of this substitution. The effect is typically, although not always, a loss of efficiency, at least when there are  $m \leq 4$  replicates. The clearest results occur when the variance does not depend on the mean. Normal theory maximum likelihood is a weighted regression of squared residuals; the corresponding method would be a weighted regression based on sample variances. Using the latter entails a loss of efficiency, no matter what the underlying distribution. For normally distributed data, the efficiency is  $(m-1)/m$ , thus being only 50% for duplicates. For other methods, using the replicate standard deviations can be more efficient. This is particularly true of a method due to Harvey (1976), which is based on the logarithm of absolute

residuals. A small absolute residual, which seems to always occur in practice, can wreak havoc with this method. This is consistent with our influence function calculations, so that we suggest some trimming of the smallest absolute residuals before applying Harvey's method.

(IV). Our results indicate that the precision of estimates  $\theta$  is approximately independent of  $\sigma$ . Also, in the power of the mean model (1.2), the efficiency of a regression estimator improves as the relative range of values of the mean response increases; efficiency depends on the spread of the logarithms of means, not their actual values. This helps explain why in assays, estimating variances is typically much harder than estimating means.

In Section 2 we discuss the prediction and calibration problems as a motivating example of a situation in which variance function estimation is of key importance. In Section 3 we describe a number of methods for estimation of  $\theta$ . We do not discuss robust methods, see Giltinan, Carroll and Ruppert (1986). In Section 4 we present an asymptotic theory for a general estimator of  $\theta$  whose construction encompasses the methods of Section 3. Section 5 contains examples of specific applications of our theory and a discussion of the implications of our formulation. Sketches of proofs are presented in Appendix A.

## 2. AN EXAMPLE: THE ROLE OF VARIANCE ESTIMATION IN PREDICTION AND CALIBRATION PROBLEMS

One example in which heterogeneity of variation occurs is in calibration experiments in the physical and biological sciences, in which one fits a model such as (1.1) to a sample  $\{y_i, x_i\}$ ,  $i = 1, \dots, N$ . The  $\{x_i\}$  may be concentrations of a substance and the  $\{y_i\}$  counts or intensity levels which vary with concentration. The interest lies in using the estimated regression to make

inference about a pair  $\{y_0, x_0\}$  which is independent of the original data set. One may wish to obtain point and interval predictors for  $y_0$  in the case  $x_0$  is known (prediction) or estimate  $x_0$  in the case  $y_0$  only is known (calibration), see Rosenblatt and Speigelman (1981). As a motivating example for considering estimation of variance functions as an independent problem, we describe the primary role of form and estimation of the variance function in construction of prediction/calibration intervals in the case of heteroscedasticity.

Throughout this discussion assume  $x_i \equiv z_i$  so that we may write the variance function as  $g(x_i, \beta, \theta)$ , and assume that the data are approximately normally distributed. Given  $x_0$ , the standard point estimate of the response  $y_0$  is  $f(x_0, \hat{\beta})$ , where  $\hat{\beta}$  is some estimate for  $\beta$ . For any consistent estimator  $\hat{\beta}$  of  $\beta$ , under (1.1) the variance in the error made by the prediction is, for large sample sizes,  $\text{var}\{y_0 - f(x_0, \hat{\beta})\} \approx \sigma^2 g^2(x_0, \beta, \theta)$ , so that the error in prediction is determined mainly by the variance function  $\sigma^2 g^2(x_0, \beta, \theta)$  and not the original data set itself. An approximate  $(1-\alpha)100\%$  confidence interval for  $y_0$  is  $I(x_0) = \{ \text{all } y \text{ in the interval } f(x_0, \hat{\beta}) \pm t_{1-\alpha/2}^{N-p} \hat{\sigma} g(x_0, \hat{\beta}, \hat{\theta}) \}$ , here  $t_{1-\alpha/2}^{N-p}$  is the  $(1-\alpha/2)$  percentage point of the  $t$  distribution with  $(N-p)$  degrees of freedom, and  $\hat{\sigma}$  and  $\hat{\theta}$  are estimates. If the parameters are estimated by a weighted analysis such as generalized least squares assuming (1.1), all estimates are consistent and the prediction interval becomes

$$(2.1) \quad I(x_0) \approx \{ \text{all } y \text{ in the interval } f(x_0, \beta) \pm t_{1-\alpha/2}^{N-p} \sigma g(x_0, \beta, \theta) \}.$$

If one were to ignore the heterogeneity, the interval would be given by  $I_U(x_0) = \{ \text{all } y \text{ in the interval } f(x_0, \hat{\beta}) \pm t_{1-\alpha/2}^{N-p} \hat{\sigma} \}$ . For an unweighted analysis, however,  $\sigma^2$  would be estimated by the unweighted mean squared error  $\hat{\sigma}_U^2 \approx \sigma^2 N^{-1} \sum g^2(x_i, \beta, \theta) = \sigma^2 g_N^2$  for large  $N$ . Thus, the unweighted prediction interval satisfies

Figure 1.

Approximate form of prediction intervals for a  
linear mean response function based on  
unweighted (ignoring heteroscedasticity) and  
weighted (as in (1.1)) regression fits.



UNWEIGHTED = DASHED, WEIGHTED = SOLID

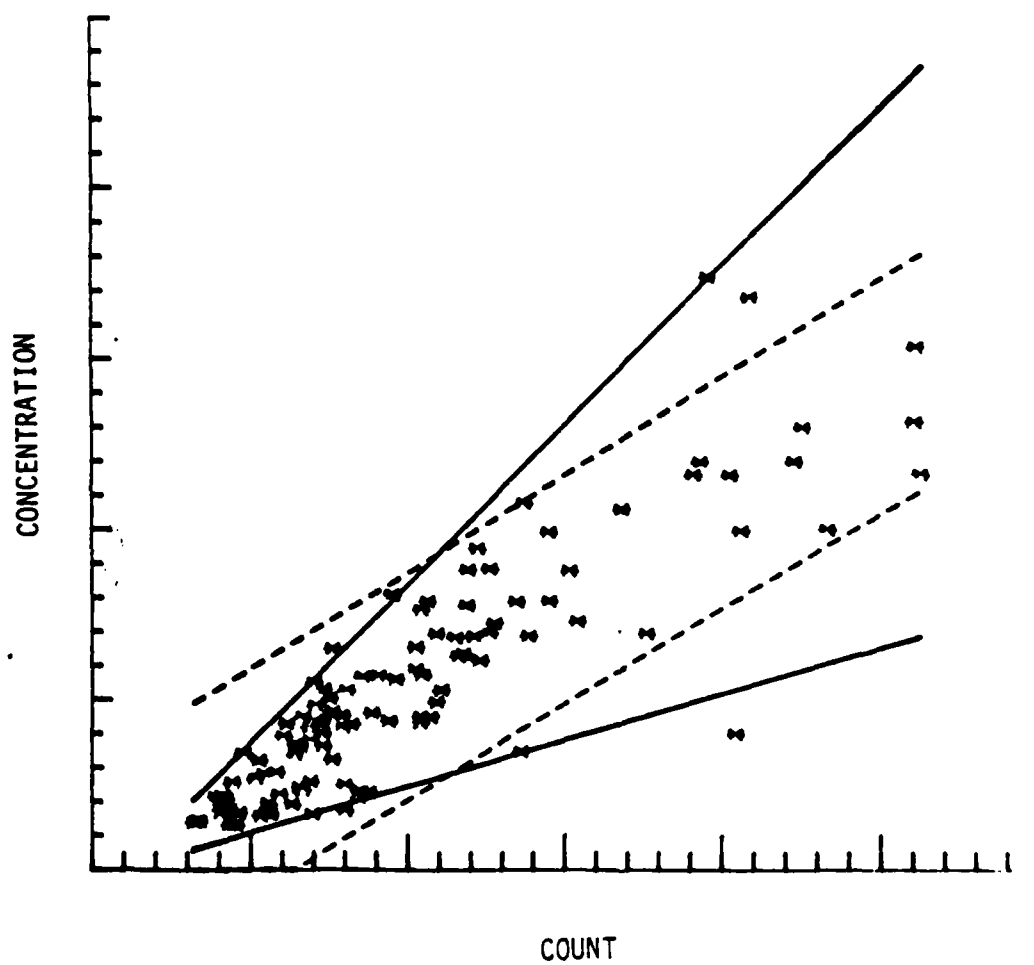
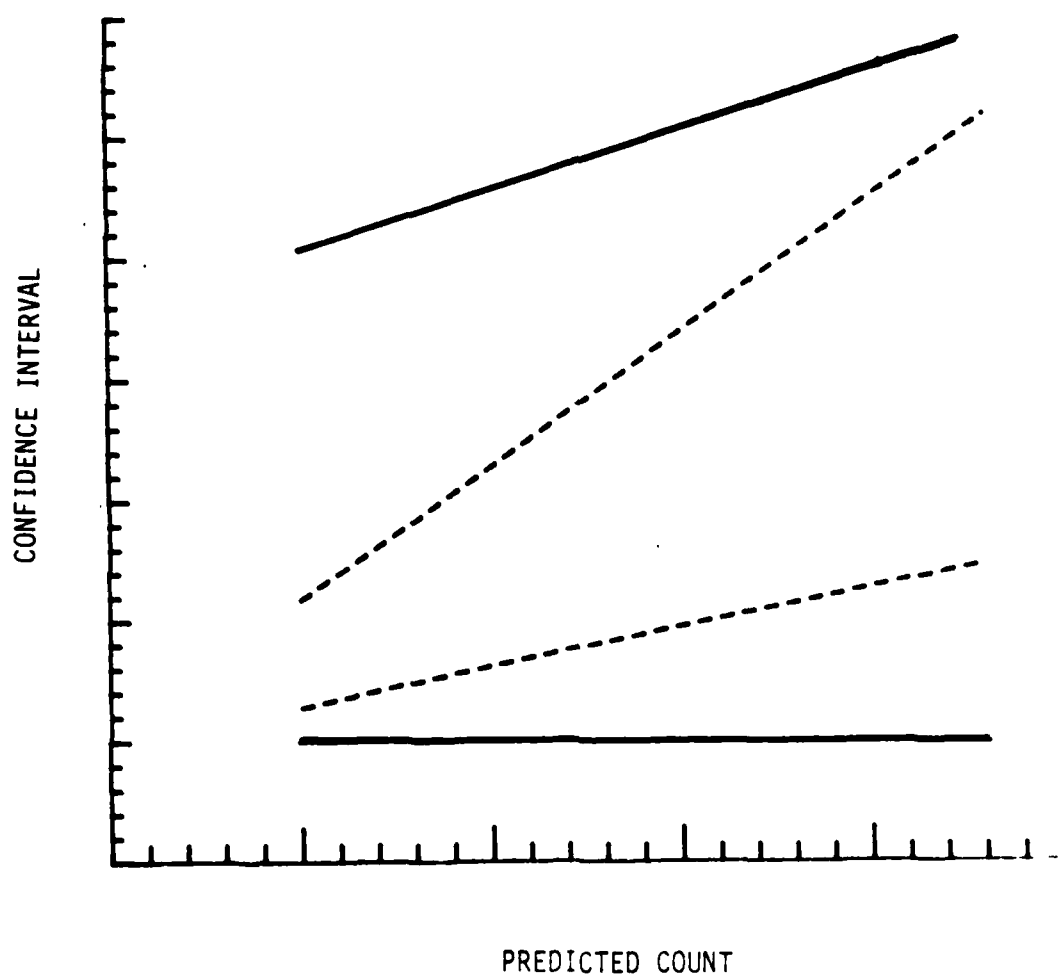


Figure 2.

Approximate form of calibration confidence intervals  
for a linear mean response function based on  
unweighted (ignoring heteroscedasticity) and  
weighted (as in (1.1)) regression fits.

UNWEIGHTED = SOLID, WEIGHTED = DASHED



$$(2.2) \quad I_U(x_0) \approx \{ \text{all } y \text{ in the interval } f(x_0, \beta) \pm t_{1-\alpha/2}^{N-p} \sigma \varepsilon_N \}.$$

Comparing (2.1) and (2.2), we see that where the variability is small, the unweighted interval will be too long and hence pessimistic, and conversely where the variance is large. Figure 1 illustrates a typical case of this phenomenon for the simple situation of an approximately linear mean response function where variability increases with mean response.

The situation is the same for calibration. For simplicity in discussing calibration, assume  $f(x, \beta)$  is strictly increasing or decreasing in  $x$ . Given  $y_0$ , the usual estimate of  $x_0$  is that value satisfying  $y_0 = f(x, \hat{\beta})$ . The common confidence interval for  $x_0$  is the set of all  $x$  values for which  $y_0$  falls in the prediction interval  $I(x)$ ; this interval is actually a  $(1-\alpha)100\%$  confidence interval for the unknown  $x_0$ . Again, the effect of not weighting is intervals which are too long for  $x_0$  where the variance is small and the opposite when the variance is large. We are not familiar with any extensive investigation of calibration confidence intervals for heteroscedastic models, although see Watters, Carroll and Spiegelman (1987). Figure 2 represents an example of this phenomenon in the situation of Figure 1 for the region of small variance.

The key point of this discussion is that when heterogeneity of variance is present, how well one models and estimates the variances will have substantial impact on prediction and calibration based on the estimated mean response, since the form of the intervals depends on the form of the variance function. Some theoretical work has been done verifying the implications of this discussion; for an investigation of how the statistical properties of estimators for calibration quantities depend on those of the estimator  $\theta$ , see Carroll, Davidian and Smith (1986) and Carroll (1987).

### 3. ESTIMATION OF $\theta$

We now discuss the form and motivation for several estimators of  $\theta$  in (1.1). In what follows, let  $\hat{\beta}_*$  be a preliminary estimator for  $\beta$ . This could be unweighted least squares or the current estimate in an iterative reweighted least squares calculation. Let  $\epsilon_i = \{Y_i - f(x_i, \beta)\} / \{g(z_i, \beta, \theta)\}$  denote the errors so that  $E \epsilon_i = 0$  and  $E \epsilon_i^2 = 1$ , and denote the residuals by  $r_i = Y_i - f(x_i, \hat{\beta}_*)$ . We consider some methods requiring  $m_i \geq 2$  replicates at each of  $M$  design points; for simplicity, we consider only the case of equal replication  $m_i \equiv m$  and write in obvious fashion  $\{Y_{ij}\}$ ,  $j = 1, \dots, m$ , to denote the  $m$  observations at  $x_i$  where appropriate, so that  $N = Mm$  is the total number of observations. In this case, let  $\bar{Y}_i$  and  $s_i$  denote the sample mean and standard deviation at  $x_i$ . For consistency of exposition, however, we denote the sum over all observations as

$$\sum_{i=1}^N \text{ instead of } \sum_{i=1}^M \sum_{j=1}^m.$$

When we speak of replacing absolute residuals  $\{|r_i|\}$  by sample deviations  $\{s_i\}$  in the case of replication,  $|r_i|$  or  $s_i$  appears  $m$  times in the sum.

#### 3.1 Regression Methods

Pseudo-likelihood. Given  $\hat{\beta}_*$ , the pseudo-likelihood estimator maximizes the normal log-likelihood  $\ell(\hat{\beta}_*, \theta, \sigma)$ , where

$$(3.1) \quad \begin{aligned} \ell(\beta, \theta, \sigma) = & -N \log \sigma - \sum_{i=1}^N \log \{g(z_i, \beta, \theta)\} \\ & - (2\sigma^2)^{-1} \sum_{i=1}^N \{Y_i - f(x_i, \beta)\}^2 / g^2(z_i, \beta, \theta). \end{aligned}$$

see Carroll and Ruppert (1982a). Here the term "pseudo-likelihood" is used as

in Gong and Samaniego (1981). Generalizations of pseudo-likelihood for robust estimation have been studied by Carroll and Ruppert (1982a) and Giltinan, Carroll and Ruppert (1986).

Least squares on squared residuals. Besides pseudo-likelihood, other methods using squared residuals have been proposed. The motivation for these methods is that the squared residuals have approximate expectation  $\sigma^2 g^2(z_i, \beta, \theta)$ , see Jobson and Fuller (1980) and Amemiya (1977). This suggests a nonlinear regression problem in which the "responses" are  $\{r_i^2\}$  and the "regression function" is  $\sigma^2 g^2(z_i, \hat{\beta}_*, \theta)$ . The estimator  $\hat{\theta}_{SR}$  minimizes in  $\theta$  and  $\sigma$

$$\sum_{i=1}^N \{r_i^2 - \sigma^2 g^2(z_i, \hat{\beta}_*, \theta)\}^2.$$

For normal data the squared residuals have approximate variance  $\sigma^4 g^4(z_i, \beta, \theta)$ ; in the spirit of generalized least squares, this suggests the weighted estimator which minimizes in  $\theta$  and  $\sigma$

$$(3.2) \quad \sum_{i=1}^N \{r_i^2 - \sigma^2 g^2(z_i, \hat{\beta}_*, \theta)\}^2 / g^4(z_i, \hat{\beta}_*, \hat{\theta}_*),$$

where  $\hat{\theta}_*$  is a preliminary estimator for  $\theta$ ,  $\hat{\theta}_{SR}$  for example. Full iteration, when it converges, would be equivalent to pseudo-likelihood.

Accounting for the effect of leverage. One objection to methods such as pseudo-likelihood and least squares based on squared residuals is that no compensation is made for the loss of degrees of freedom associated with preliminary estimation of  $\beta$ . For example, the effect of applying pseudo-likelihood directly seems to be a bias depending on  $p/N$ . For settings such as fractional factorials where  $p$  is large relative to  $N$  this bias could be substantial.

Bayesian ideas have been used to account for loss of degrees of freedom;

see Harville (1977) and Patterson and Thompson (1974). When  $g$  does not depend on  $\beta$ , the restricted maximum likelihood approach of the latter authors suggests in our setting one estimate  $\theta$  from the mode of the marginal posterior density for  $\theta$  assuming normal data and a prior for the parameters proportional to  $\sigma^{-1}$ . When  $g$  depends on  $\beta$ , one may extend the Bayesian arguments and use a linear approximation as in Box and Hill (1974) and Beal and Sheiner (1986) to define a restricted maximum likelihood estimator.

Let  $Q$  be the  $N \times p$  matrix with  $i$ th row  $f_{\beta}(x_i, \beta)^t / g(z_i, \beta, \theta)$ , where  $f_{\beta}(x_i, \beta) = \partial / \partial \beta \{f(x_i, \beta)\}$ , and let  $H = Q(Q^t Q)^{-1} Q^t$  be the "hat" matrix with diagonal element  $h_{ii} = h_{ii}(\beta, \theta)$ ; the values  $\{h_{ii}\}$  are the leverage values. It turns out that the restricted maximum likelihood estimator is equivalent to an estimator obtained by modifying pseudo-likelihood to account for the effect of leverage. This characterization, while not unexpected, is new; we derive this estimator and its equivalence to a modification of pseudo-likelihood in Appendix B.

The least squares approach using squared residuals can also be modified to show the effect of leverage. Jobson and Fuller (1980) essentially note that for nearly normally distributed data we have the approximations

$$\begin{aligned} E r_i^2 &\approx \sigma^2 (1 - h_{ii}) g^2(z_i, \beta, \theta), \\ \text{var } r_i^2 &\approx 2 \sigma^4 (1 - h_{ii})^2 g^4(z_i, \beta, \theta). \end{aligned}$$

To exploit these approximations modify (3.2) to minimize in  $\theta$  and  $\sigma$

$$(3.3) \quad \sum_{i=1}^N \{r_i^2 - \sigma^2 (1 - \hat{h}_{ii}) g^2(z_i, \hat{\beta}_{*}, \theta)\}^2 / \{(1 - \hat{h}_{ii})^2 g^4(z_i, \hat{\beta}_{*}, \hat{\theta}_{*})\},$$

where  $\hat{h}_{ii} = h_{ii}(\hat{\beta}_{*}, \hat{\theta}_{*})$  and  $\hat{\theta}_{*}$  is a preliminary estimator for  $\theta$ . An asymptotically equivalent variation of this estimator in which one sets the derivatives of (3.3) with respect to  $\theta$  and  $\sigma$  equal to 0 and then replaces  $\hat{\theta}_{*}$  by

$\theta$  can be seen to be equivalent to pseudo-likelihood in which one replaces standardized residuals by studentized residuals. While this estimator also takes into account the effect of leverage, it is different from restricted maximum likelihood.

Least squares on absolute residuals. Squared residuals are skewed and long-tailed, which has lead many authors to propose using absolute residuals to estimate  $\theta$ ; see Glejser (1969) and Theil (1971). Assume that

$$E|Y_i - f(x_i, \beta)| = \eta g(z_i, \beta, \theta),$$

which is satisfied if the errors  $\{\epsilon_i\}$  are independent and identically distributed. Mimicking the least squares approach based on squared residuals, one obtains the estimator  $\hat{\theta}_{AR}$  by minimizing in  $\eta$  and  $\theta$

$$\sum_{i=1}^N \{|r_i| - \eta g(z_i, \hat{\beta}_*, \theta)\}^2.$$

In analogy to (3.2), the weighted version is obtained by minimizing

$$\sum_{i=1}^N \{|r_i| - \eta g(z_i, \hat{\beta}_*, \theta)\}^2 / g^2(z_i, \hat{\beta}_*, \hat{\theta}_*).$$

where  $\hat{\theta}_*$  is a preliminary estimator for  $\theta$ , probably  $\hat{\theta}_{AR}$ . As for least squares estimation based on squared residuals, one presumably could modify this approach to account for the effect of leverage.

Logarithm method. The suggestion of Harvey (1976) is to exploit the fact that the logarithm of the absolute residuals has approximate expectation  $\log \{g(z_i, \beta, \theta)\}$ . Estimate  $\theta$  by ordinary least squares regression of  $\log |r_i|$  on  $\log \{g(z_i, \hat{\beta}_*, \theta)\}$ , since if the errors are independent and identically distributed, the regression should be approximately homoscedastic. If one of



the residuals is near zero the regression could be adversely affected by a large "outlier," hence in practice one might wish to delete a few of the smallest absolute residuals, perhaps trimming the smallest few percent.

### 3.2 Other methods

Besides squares and logarithms of absolute residuals, other transformations could be used. For example, the square root and  $2/3$  root would typically be more normally distributed than the absolute residuals themselves. Such transformations appear to be useful, although they have not been used much to our knowledge. Our asymptotic theory applies to such transformations.

In a parametric model such as (1.1), joint maximum likelihood estimation is possible, where we use the term maximum likelihood to mean normal theory maximum likelihood. When the variance function does not depend on  $\beta$ , it can be easily shown that maximum likelihood is asymptotically equivalent to weighted least squares methods based on squared residuals. In the situation in which the variance function depends on  $\beta$  this is not the case. In this setting, it has been observed by Carroll and Ruppert (1982b) and McCullagh (1983) that while maximum likelihood estimators enjoy asymptotic optimality when the model and distributional assumptions are correct, the maximum likelihood estimator of  $\beta$  can suffer problems under departures from these assumptions. This suggests that joint maximum likelihood estimation should not be applied blindly in practice. The theory of the next section shows the asymptotic equivalence of maximum likelihood with other methods in a simplifying special case. Based on this theory, we tend to prefer weighted regression methods even when the data are approximately normal for reasons of relative computational simplicity.

While we have chosen to describe the methods of Section 3.1 as "regression methods," asymptotically equivalent versions of such methods may be derived by

considering maximum likelihood assuming some underlying distribution. For example, the form of the weighted squared residuals method is that of normal theory maximum likelihood with  $\beta$  known and  $\hat{\theta}_*$  replaced by  $\theta$  (pseudo-likelihood); the form of the weighted absolute residual method is that of maximum likelihood assuming  $\beta$  known and  $\hat{\theta}_*$  replaced by  $\theta$  under the double exponential distribution. Thus, what we term a regression method may be viewed as an approximation to maximum likelihood assuming a particular distribution. We feel that the regression interpretation is a much more appealing and natural motivation, since no particular distribution need be considered to obtain the form of the estimators, only the mean-variance relationship.

Another joint estimation method is the extended quasi-likelihood of Nelder and Pregibon (1987) also described in McCullagh and Nelder (1983). This estimator is based on assuming a class of distributions "nearly" containing skewed distributions such as the Poisson and gamma. While it may be viewed as iteration between estimation of  $\theta$  and  $\sigma$  and generalized least squares for  $\beta$ , technically this scheme does not fit in the general framework of the next section; an asymptotic theory has been developed elsewhere, see Davidian and Carroll (1987). A related formulation is given by Efron (1986).

Methods requiring replicates at each design point have been proposed in the assay literature. These methods do not depend on the postulated form of the regression function; one reason that this may be advantageous is that in many assays along with observed pairs  $(Y_{ij}, x_i)$  there will also be pairs in which only  $Y_{ij}$  is observed. A popular and widely used method is that of Rodbard and Frazier (1975). If we assume

$$(3.4) \quad g(z_i, \beta, \theta) = g(\mu_i, z_i, \theta),$$

as in, for example, (1.2) or (1.3), the method is identical to the logarithm

method previously discussed except that one replaces  $|r_i|$  by the sample standard deviation  $s_i$  and  $f(x_i, \hat{\beta}_*)$  in the "regression" function by the sample mean  $\bar{Y}_i$ . As a motivation for this and the method of Harvey, consider that under (1.2)  $\theta$  is simply the slope parameter for a simple linear regression.

As an alternative, under the assumption of independence and (3.4), the modified maximum likelihood method of Raab (1981) estimates  $\theta$  by joint maximization in the  $(M+r+1)$  parameters  $\sigma^2, \theta, \mu_1, \dots, \mu_M$  of the "modified" normal likelihood

$$(3.5) \quad \prod_{i=1}^M \{2\pi\sigma^2 g^2(\mu_i, z_i, \theta)\}^{(m-1)/2} \exp[-\sum_{j=1}^m (Y_{ij} - \mu_i)^2 / \{2\sigma^2 g^2(\mu_i, z_i, \theta)\}]$$

The modification serves to make the estimator of  $\sigma$  unbiased. The idea here is to improve upon the regression method of Rodbard by appealing to a maximum likelihood approach which, despite a parameter space increasing as the number of design points, is postulated to have reasonable properties. A related method is that in which  $\theta$  and  $\sigma$  are estimated by maximizing (3.5) with  $\mu_i$  replaced by  $\bar{Y}_i$ , the motivation being computational ease and evidence that this estimator may not be too different from that of Raab in practice, see Sadler and Smith (1985).

Table 1 contains a summary of some of the common methods for variance function estimation and their formulations.

#### 4. AN ASYMPTOTIC THEORY OF VARIANCE FUNCTION ESTIMATION

In this section we construct an asymptotic theory for a general class of regression-type estimators for  $\theta$ . Since our major interest lies in obtaining general insights, we do not state technical assumptions or details. In what

follows, in the case of replication  $N \rightarrow \infty$  in such a way that  $m$  remains fixed. The reader uninterested in this development may wish to skip to Section 5, where conclusions and implications of the theory are presented.

#### 4.1 Methods based on transformations of absolute residuals

Write  $d_i(\beta) = |Y_i - f(x_i, \beta)|$ . Let  $T$  be a smooth function and define  $M_i$  by

$$M_i(\eta, \theta, \beta) = E [ T\{d_i(\beta)\} ],$$

where  $\eta$  is a scale parameter which is usually a function of  $\sigma$  only. We consider estimation of the more general parameter  $\eta$  instead of  $\sigma$  itself for ease of exposition, and since  $\sigma$  is estimated jointly with  $\theta$  in regression models, our theory focuses on expansions for  $\eta$  and  $\theta$  jointly. If  $\hat{\eta}_*$ ,  $\hat{\theta}_*$  and  $\hat{\beta}_*$  are any preliminary estimators for  $\eta$ ,  $\theta$ , and  $\beta$ , define  $\hat{\eta}$  and  $\hat{\theta}$  to be the solutions of

$$(4.1) \quad N^{-1/2} \sum_{i=1}^N H_i(\eta, \theta, \hat{\beta}_*) \{T\{d_i(\hat{\beta}_*)\} - M_i(\eta, \theta, \hat{\beta}_*)\} / V_i(\hat{\eta}, \hat{\theta}, \hat{\beta}_*) = 0,$$

where  $V_i(\eta, \theta, \beta)$  is a smooth function and  $H_i$  is a smooth function which for the estimators of Section 3 is the partial derivative of  $M_i$  with respect to  $(\eta, \theta)$ . In what follows, we suppress the arguments of the functions  $M_i$ ,  $V_i$ , etc. when they are evaluated at the true values  $\eta$ ,  $\theta$ , and  $\beta$ . Specific examples are considered in the next section.

The class of estimators solving (4.1) includes directly or includes an asymptotically equivalent version of the estimators of Section 3.1. For methods which account for the effect of leverage,  $M_i$ ,  $V_i$  and  $H_i$  will depend on the  $h_{ii}$ . In this case we need the additional assumption that if  $h = \max \{h_{ii}\}$ , then  $N^{1/2}h$  converges to zero.

**Theorem 4.1.** Let  $\hat{\eta}_*$ ,  $\hat{\theta}_*$  and  $\hat{\beta}_*$  be  $N^{1/2}$  consistent for estimating  $\eta$ ,  $\theta$  and  $\beta$ . Let  $\dot{T}$  be the derivative of  $T$  and define

$$\begin{aligned} C_i &= H_i [T\{d_i(\beta)\} - M_i] / V_i; \\ B_{1,N} &= N^{-1} \sum_{i=1}^N H_i H_i^t / V_i; \\ B_{2,N} &= -N^{-1} \sum_{i=1}^N (H_i / V_i) \partial / \partial \beta \{M_i(\eta, \theta, \beta)\}; \\ B_{3,N} &= -N^{-1} \sum_{i=1}^N (H_i / V_i) f_{\beta}(x_i, \beta) E [\dot{T}\{d_i(\beta)\} \text{sign}(\epsilon_i)]. \end{aligned}$$

Then, under regularity conditions as  $N \rightarrow \infty$ ,

$$(4.2) \quad B_{1,N} N^{1/2} \begin{bmatrix} \hat{\eta} - \eta \\ \hat{\theta} - \theta \end{bmatrix} = N^{-1/2} \sum_{i=1}^N C_i + (B_{2,N} + B_{3,N}) N^{1/2} (\hat{\beta}_* - \beta) + o_p(1). \quad \square$$

We may immediately make some general observations about the estimator  $\hat{\theta}$  solving (4.1). Note that if the variance function does not depend on  $\beta$ , then  $M_i$  does not depend on  $\beta$  and hence  $B_{2,N} \equiv 0$ . For the estimators of Section 2.1,  $\dot{T}$  is an odd function. Thus, if the errors  $\{\epsilon_i\}$  are symmetrically distributed,  $E[\dot{H}_i\{d_i(\beta)\} \text{sign}(\epsilon_i)] = 0$  and hence  $B_{3,N} \equiv 0$ .

**Corollary 4.1(a).** Suppose that the variance function does not depend on  $\beta$  and the errors are symmetrically distributed. Then the asymptotic distributions of the regression estimators of Section 3.1 do not depend on the method used to obtain  $\hat{\beta}_*$ . If both of these conditions do not hold simultaneously, then the asymptotic distributions will depend in general on the method of estimating  $\beta$ .

□

The implication is that in the situation for which the variance function does not depend on  $\beta$  and the data are approximately symmetrically distributed, for large sample sizes the preliminary estimator for  $\beta$  will play little role in

determining the properties of  $\hat{\theta}$ . Note also from (4.2) that for weighted methods, the effect of the preliminary estimator of  $\theta$  is asymptotically negligible regardless of the underlying distributions.

The preliminary estimator  $\hat{\beta}_*$  might be the unweighted least squares estimator, a generalized least squares estimator or some robust estimator. See, for example, Huber (1981) and Giltinan, Carroll and Ruppert (1986) for examples of robust estimators for  $\beta$ . For some vectors  $\{v_{N,i}\}$ , these estimators admit an asymptotic expansion of the form

$$(4.3) \quad N^{1/2}(\hat{\beta}_* - \beta) = N^{-1/2} \sum_{i=1}^N \Psi(v_{N,i}, \epsilon_i) + o_p(1).$$

Here  $\Psi$  is odd in the argument  $\epsilon$ . In case the variance function depends on  $\beta$ ,  $B_{2,N} \neq 0$  in general; however, if the errors are symmetrically distributed and  $\hat{\beta}_*$  has expansion of form (4.3), then the two terms on the right-hand side of (4.2) are asymptotically independent. The following is then immediate.

Corollary 4.1(b). Suppose that the errors are symmetrically distributed and that  $\hat{\beta}_*$  has an asymptotic expansion of the form (4.3). Then for the estimators of Section 3.1, the asymptotic covariance matrix of  $\hat{\theta}$  is a monotone nondecreasing function of the asymptotic covariance matrix of  $\hat{\beta}_*$ .  $\square$

By the Gauss-Markov theorem and the results of Jobson and Fuller (1980) and Carroll and Ruppert (1982a), the implication of Corollary 4.1(b) is that using unweighted least squares estimates of  $\beta$  will result in inefficient estimates of  $\theta$ . This phenomenon is exhibited in small samples in a Monte Carlo study of Carroll, Davidian and Smith (1986). If one starts from the unweighted least squares estimate, one ought to iterate the process of estimating  $\theta$  -- use the current value  $\hat{\beta}_*$  to estimate  $\theta$  from (4.1), use these  $\hat{\beta}_*$  and  $\hat{\theta}$  to obtain an

updated  $\hat{\beta}_*$  by generalized least squares and repeat the process  $\ell - 1$  more times. It is clear that the asymptotic distribution of  $\hat{\theta}$  will be the same for  $\ell \geq 2$  with larger asymptotic covariance for  $\ell = 1$ , so in principle one ought to iterate this process at least twice. See Carroll, Wu and Ruppert (1987) for more on iterating generalized least squares.

#### 4.2 Methods based on sample standard deviations

Assume replication, and as before let  $\{s_i\}$  be the sample standard deviations at each  $x_i$ , which themselves have been proposed as estimators of the variance in generalized least squares estimation of  $\beta$ . This can be disastrous, see Jacquez, Mather and Crawford (1968). When replication exists, however, practitioners feel comfortable with the notion that the  $\{s_i\}$  may be used as a basis for estimating variances; thus, one might reasonably seek to estimate  $\theta$  by replacing  $d_i(\hat{\beta}_*)$  by  $s_i$  in (4.1).

The following result is almost immediate from the proof of Theorem 4.1 in Appendix A.

Theorem 4.2. If  $d_i(\hat{\beta}_*)$  is replaced by  $s_i$  in (4.1), then under the conditions of Theorem 4.1 the resulting estimator for  $\theta$  satisfies (4.2) with  $B_{3,N} \equiv 0$  and the redefinitions

$$(4.4a) \quad C_i = (H_i/V_i)\{T(s_i) - M_i\};$$

$$(4.4b) \quad M_i = E\{T(s_i)\} = M_i(\eta, \theta, \beta).$$

□

If the errors are symmetrically distributed, then from (4.2) and Theorem 4.2, whether one is better off using absolute residuals or sample standard deviations in the methods of Section 3.1 depends only on the differences

between the expected values and variances of  $T\{d_i(\beta)\}$  and  $T(s_i)$ . In Section 5 we exhibit such comparisons explicitly and show that absolute residuals can be preferred to sample standard deviations in situations of practical importance.

#### 4.3 Methods not depending on the regression function

We assume throughout this discussion that the variance function has form (3.4) and replication is available. From Section 3.1 we see that the "regression function" part of the estimating equations depends on  $f(x_i, \hat{\beta}_*)$ , so that in the general equation (4.1)  $M_i$ ,  $V_i$  and  $H_i$  all depend on  $f(x_i, \hat{\beta}_*)$ . In some settings, one may not postulate a form for the  $\mu_i$  for estimating  $\theta$ ; the method of Rodbard and Frazier (1975), for example, uses  $s_i$  in place of  $d_i(\hat{\beta}_*)$  as in Section 4.2 and replaces  $f(x_i, \hat{\beta}_*)$  by the sample mean  $\bar{Y}_i$ . We now consider the effect of replacing predicted values by sample means for the general class (4.1).

The presence of the sample means in the variance function in (4.1) requires more complicated and restrictive assumptions than the usual large sample asymptotics applied heretofore. The method of Rodbard and Frazier and the general method (4.1) with sample means are functional nonlinear errors in variables problems as studied by Wolter and Fuller (1982) and Stefanski and Carroll (1985). Standard asymptotics for these problems correspond to letting  $\sigma$  go to zero at rate  $N^{-1/2}$ . In Section 4.4 we discuss the practical implications of  $\sigma$  being small; for now, we state the following result.

Theorem 4.3. Suppose that we replace  $f(x_i, \hat{\beta}_*)$  by  $\bar{Y}_i$  in  $M_i$ ,  $V_i$  and  $H_i$  in (4.1) and adopt the assumptions of Theorems 4.1 and 4.2. Further, suppose that as  $N \rightarrow \infty$ ,  $\sigma \rightarrow 0$  simultaneously and

$$(i) \quad N^{1/2}\sigma \rightarrow \lambda, \quad 0 \leq \lambda < \infty;$$



- (ii)  $N^{1/2} \sum_{i=1}^N C_i$  has a nontrivial asymptotic normal limit distribution;
- (iii) The  $\{\epsilon_i\}$  are symmetric and i.i.d.;
- (iv)  $\{|\bar{Y}_{i.} - \mu_1| / \sigma\}^2$  has uniformly bounded  $k$  moments, some  $k > 2$ .

Then the results of Theorems 4.1 and 4.2 hold with  $B_{2,N} = B_{3,N} \equiv 0$ .  $\square$

This result shows that under certain restrictive assumptions, one may replace predicted values by sample means under replication; however, it is important to realize that the assumption of small  $\sigma$  is not generally valid and hence the use of sample means may be disadvantageous in situations where these asymptotics do not apply. Further, relaxation of assumption (iii) will result in an asymptotic bias in the asymptotic distribution of the estimator not present for estimators based on residuals regardless of the assumption of symmetry; see Appendix A.

The estimator of Raab (1981) discussed in Section 3.2 is also a functional nonlinear errors in variables estimator, complicated by a parameter space with size of order  $N$ . Sadler and Smith (1985) have observed that the Raab estimator is often indistinguishable from their estimator with  $\mu_1$  replaced by  $\bar{Y}_{i.}$  in (3.5); such an estimator is contained in the general class (4.1). Davidian (1986) has shown that under the asymptotics of Theorem 4.3 and additional regularity conditions that the two estimators are asymptotically equivalent in an important special case. We may thus consider the result of Theorem 4.3 relevant to this estimator.

#### 4.4 Small $\sigma$ asymptotics

In Section 4.3 technical considerations forced us to pursue an asymptotic theory in which  $\sigma$  is small. It turns out that in some situations of practical

importance these asymptotics are relevant. In particular, in assay data we have observed values for  $\sigma$  which are quite small relative to the means. Such asymptotics are used in the study of data transformations in regression. It is thus worthwhile to consider the effect of small  $\sigma$  on the results of Sections 4.1 and 4.2 and to comment on some other implications of letting  $\sigma \rightarrow 0$ .

In the situation of Theorem 4.1, if the errors are symmetrically distributed, then for the estimators of Section 3.1, if  $\sigma \rightarrow 0$  as  $N \rightarrow \infty$ , then there is no effect for estimating the regression parameter  $\beta$ . In the situation of Theorem 4.2, the errors need not even be symmetrically distributed. The major insight provided by these results is that in certain practical situations in which  $\sigma$  is small, the choice of  $\hat{\beta}_*$  may not be too important even if the variance function depends on  $\beta$ .

Small  $\sigma$  asymptotics may be used also to provide insight into the behavior of other estimators for  $\theta$  which do not fit into the general framework of (4.1). It can be shown that the extended quasi-likelihood estimator need not necessarily be consistent for fixed  $\sigma$ , but if one adopts the asymptotics of the previous section, this estimator is asymptotically equivalent to regression estimators based on squared residuals as long as the errors are symmetrically distributed. Otherwise, an asymptotic bias results which may have implications for inference for  $\theta$ . For discussion see Davidian and Carroll (1987).

The small  $\sigma$  assumption also provides an illustration of the relationship between variance function estimation and data transformations. Let  $\ell(y, \lambda) = (y^\lambda - 1)/\lambda$ , and consider the model

$$(4.5) \quad E\{ \ell(Y_1, \lambda) \} = \ell( f(x_1, \beta), \lambda ); \quad \text{var}\{ \ell(Y_1, \lambda) \} = \sigma;$$

such "transform both sides" models are proposed and motivated by Carroll and Ruppert (1984). For  $\sigma \approx 0$ ,  $E(Y_1) \approx f(x_1, \beta)$  and  $\text{var}(Y_1) \approx \sigma f(x_1, \beta)^{(1-\lambda)}$ , so

that in (1.2) we have  $\theta \approx 1 - \lambda$ . Thus, when the small  $\sigma$  assumption is relevant, (4.5) and (1.1), (1.2) represent approximately the same model.

## 5. APPLICATIONS AND FURTHER RESULTS

In Section 4 we constructed an asymptotic theory for and stated some general characteristics of regression-type estimators of  $\theta$ . In this section we use the theory to exhibit the specific forms for the various estimators of Section 3 and compare and contrast their properties. In our investigation we rely on the simplifying assumptions implied by the theory of Section 4, in particular the small  $\sigma$  asymptotic approach in which  $\sigma \rightarrow 0$  as  $N \rightarrow \infty$ . Throughout, define  $v(i, \beta, \theta) = \log g(z_i, \beta, \theta)$ , let  $v_\theta(i, \beta, \theta)$  be the column vector of partial derivatives of  $v$  with respect to  $\theta$ , let  $\xi(\beta, \theta)$  be the covariance matrix of  $v_\theta(i, \beta, \theta)$ , and let  $\tau(i, \beta, \theta) = \{1, v_\theta^t(i, \beta, \theta)\}^t$ . For simplicity, assume that the errors  $\{\epsilon_i\}$  are independent and identically distributed with kurtosis  $\kappa$ ;  $\kappa = 0$  for normality.

### 5.1 Maximum likelihood, pseudo-likelihood, restricted maximum likelihood and weighted squared residuals.

Writing  $\eta = \log \sigma$ , we have  $T(x) = x^2$ ,  $M_i = \exp(2\eta) g^2(z_i, \beta, \theta)$ ,  $V_i = M_i^2$ ,  $H_i^t = \partial M_i / \partial (\eta, \theta^t)^t$  and  $E [\dot{T}\{d_i(\beta)\} \text{sign}(\epsilon_i)] = 2 E [Y_i - f(x_i, \beta)] = 0$  so that  $B_{3,N} \equiv 0$  regardless of the underlying distributions. If  $h \rightarrow 0$  such that  $N^{1/2}h \rightarrow 0$  for methods accounting for the effect of leverage, then all of these methods admit an expansion of the form (4.2) with  $B_{3,N} = 0$ . The expansion will be different depending on whether  $\hat{\beta}_*$  is a generalized least squares estimator for  $\beta$  or full maximum likelihood, since the maximum likelihood estimator has an

expansion quadratic in the errors while that of the generalized least squares estimator is linear in the  $\{\epsilon_i\}$ , see Carroll and Ruppert (1982b). The implication is that regression methods based on iterated weighted squared residuals and full maximum likelihood are different in general asymptotically. Regardless of the underlying distributions, for fixed  $\sigma$ , Davidian (1986) has shown that the asymptotic covariance matrix of the former methods increases without bound as a function of  $\sigma$  while that of maximum likelihood remains bounded for all  $\sigma$ . Further, a simple comparison of the two covariances reveals that under reasonable conditions maximum likelihood has smaller asymptotic covariance as long as  $\kappa \leq 2$ . While these facts may suggest a preference for full maximum likelihood even away from normality, the computational and model robustness considerations mentioned earlier may make this preference tenuous. Generalized least squares and maximum likelihood estimators for  $\beta$  both satisfy  $\hat{\beta}_* - \beta = O_p(\sigma N^{-1/2})$ , so that if  $\sigma \rightarrow 0$  or  $g$  does not depend on  $\beta$ , then  $\hat{\theta}$  is asymptotically normally distributed with mean  $\theta$  and covariance matrix

$$(5.1) \quad (2 + \kappa) \{4N \xi(\beta, \theta)\}^{-1}.$$

As mentioned in Section 3, under the small  $\sigma$  asymptotics of Theorem 3.3, the extended quasi-likelihood estimator of  $\theta$  is asymptotically equivalent to the estimators here with asymptotic covariance matrix (5.1). Thus, if  $g$  does not depend on  $\beta$  or  $\sigma \rightarrow 0$ , pseudo-likelihood, weighted squared residuals, restricted maximum likelihood, maximum likelihood and, if  $\sigma \rightarrow 0$ , extended quasi-likelihood, are all asymptotically equivalent. In addition, all of these estimators have influence functions which are linear in the squared errors, indicating substantial nonrobustness.

We may also observe that these methods are preferable to unweighted regression on squared residuals. Write (5.1) as

$$(5.2) \quad (1/2 + \kappa/4) (WV^{-1}W)^{-1},$$

where  $V$  is the  $N \times N$  diagonal matrix with elements  $V_i$  and  $W$  is the  $N \times p$  matrix with  $i^{\text{th}}$  row  $H_i^t$ . For the unweighted estimator based on squared residuals, calculations similar to those above show that the asymptotic covariance matrix when either  $g$  does not depend on  $\beta$  or  $\sigma \rightarrow 0$  is given by

$$(5.3) \quad (1/2 + \kappa/4) (W^t W)^{-1} (W^t V W) (W^t W)^{-1}.$$

The comparison between (5.2) and (5.3) is simply that of the Gauss-Markov theorem, so that (5.2) is no larger than (5.3).

## 5.2 Logarithms of absolute residuals and the effect of inliers

We do not consider deletion of the few smallest absolute residuals. Here  $T(x) = \log x$  so that  $\dot{T}(x) = x^{-1}$ . Letting  $\eta = \log \sigma$  and assuming independent and identically distributed errors we have  $M_i = \eta + v(i, \beta, \theta) + E \log |\epsilon|$ ,  $V_i \equiv 1$ , and  $H_i = \tau(i, \beta, \theta)$ . Under the assumption of symmetry of the errors, with  $g$  not depending on  $\beta$  or  $\sigma \rightarrow 0$ , tedious algebra shows that  $\hat{\theta}$  is asymptotically normally distributed with mean  $\theta$  and covariance matrix

$$(5.4) \quad \text{var} \{ \log (|\epsilon|^2) \} \{ 4N \xi(\beta, \theta) \}^{-1}.$$

The influence function for this estimator is linear in the logarithm of the absolute errors. This indicates nonrobustness more for inliers than for outliers, which at the very least is an unusual phenomenon. If the errors are not symmetric then there will be an additional effect due to estimating  $\beta$  not

present for the methods of Section 5.1, even if  $g$  does not depend on  $\beta$ .

### 5.3 Weighted Absolute Residuals

Assume that the errors are independent and identically distributed and let  $\exp(\eta) = \sigma E|\epsilon|$ . Consider the weighted estimator. We have  $T(x) = x$ ,  $\dot{T}(x) = 1$ ,  $M_1 = \exp(\eta) g(z_1, \beta, \theta)$  and  $V_1 = M_1^2$ . Thus, if the errors are symmetrically distributed and either  $g$  does not depend on  $\beta$  or  $\sigma \rightarrow 0$ ,  $\hat{\theta}$  is asymptotically normally distributed with mean  $\theta$  and covariance matrix

$$(5.5) \quad \{\delta/(1 - \delta)\} \{N \xi(\beta, \theta)\}^{-1},$$

where  $\delta = \text{var } |\epsilon|$ . The influence function for this estimator is linear in the absolute errors. By an argument similar to that at the end of Section 5.1, we may conclude that when the effect of  $\hat{\beta}_*$  is negligible one should use a weighted estimator and iterate the method.

### 5.4 General transformations

One may also consider other power transformations of absolute residuals. If  $\lambda \neq 0$  is the power of absolute residuals on which the regression is based, then define  $\eta$  by  $\exp(\lambda\eta) = \sigma^\lambda E(|\epsilon|^\lambda)$  and  $T(x) = x^\lambda$ . Then  $M_1 = \exp(\lambda\eta) g^\lambda(z_1, \beta, \theta)$ ,  $V_1 = M_1^2$ . Straightforward calculations show that if the errors are symmetric and either  $g$  does not depend on  $\beta$  or  $\sigma \rightarrow 0$ , then  $\hat{\theta}$  is asymptotically normally distributed with mean  $\theta$  and asymptotic covariance matrix

$$(5.6) \quad [\text{var} (|\epsilon|^\lambda) / \{E(|\epsilon|^\lambda)^2\}] \{ \lambda^2 N \xi(\beta, \theta) \}^{-1}.$$

with influence function linear in  $|\epsilon|^\lambda$ . For square root transformations, for example,  $\lambda = 1/2$ , and from (5.1) and (5.6), the asymptotic relative efficiency of the square root transformation relative to pseudo-likelihood under normal errors is 0.693; from (5.5), the efficiency relative to weighted absolute residuals is 0.791.

At this point it is worthwhile to mention that under the simplifying assumptions of our discussion, the precision of general regression estimators does not depend on  $\sigma$ , since a general expression such as (5.6) is independent of  $\eta$ . Thus, how well we estimate  $\theta$  in many practical cases will be approximately independent of  $\sigma$ . Furthermore, when the power of the mean model for variance (1.2) holds,  $v_\theta(1, \beta, \theta) = \log \mu_1$ , so that  $\xi(\beta, \theta)$  is the limiting variance of the  $\{\log \mu_1\}$ . From the general expression (5.6), the precision with which one can estimate  $\theta$  depends only on the relative spread of the mean responses, not their actual sizes, and clearly this spread must be fairly substantial so that the spread of the logarithms of the means will be so as well. The implications are that for (1.2), the design will play an important role in efficiency of estimation of  $\theta$ , and in some practical situations we may not be able to estimate  $\theta$  well no matter which estimator we employ.

### 5.5 Comparison of methods based on residuals

We assume that the errors are symmetric and independent and identically distributed and that either  $g$  does not depend on  $\beta$  or  $\sigma$  is small. By (5.1), (5.4) and (5.5), the asymptotic relative efficiency of the three methods depends only on the distribution of the errors. For normal errors, using absolute residuals results in a 12% loss in efficiency while for standard double exponential errors there is a 25% gain in efficiency for using absolute residuals. For normal errors, the logarithm method represents a 59% loss of

efficiency with respect to pseudo-likelihood.

In Table 2 we present asymptotic relative efficiencies for various contaminated normal distributions. The asymptotic relative efficiency of the weighted absolute residual method to pseudo-likelihood is the same as the asymptotic relative efficiency of the mean absolute deviation with respect to the sample variance for a single sample, see Huber (1981, page 3); the first column of the table is thus identical to that of Huber. The table shows that while at normality neither the absolute residuals nor the logarithm methods are efficient, a very slight fraction of "bad" observations is enough to offset the superiority of squared residuals in a dramatic fashion. For example, just two bad observations in 1000 negate the superiority of squared residuals. If 1% or 5% of the data are "bad," absolute residuals and the logarithm method, respectively, show substantial gains over squared residuals. The implication is that while it is commonly perceived that methods based on squared residuals are to be preferred in general, these methods can be highly non-robust. Our formulation includes this result for maximum likelihood, showing its inadequacy under slight departures from the assumed distributional structure. We also include asymptotic relative efficiencies for appropriately weighted residual methods based on square, cube and  $2/3$  roots to pseudo-likelihood using (5.6) and observe that these methods also exhibit comparative robustness to contamination.

#### 5.6 Methods based on sample standard deviations

Assume that  $m \geq 2$  replicate observations are available at each design point. In practice,  $m$  is usually small, see Raab (1981). We compare using absolute residuals to using sample standard deviations in the estimators of Section 3.1. We assume that one is fairly confident in the postulated form of



the model, thus viewing methods based on sample standard deviations as not taking full advantage of the information available. For simplicity, assume that the errors are independent and identically and symmetrically distributed and that either  $g$  does not depend on  $\beta$  or  $\sigma$  is small. If the errors are not symmetric and  $\sigma$  is not small or the variance depends on  $\beta$ , using sample standard deviations presumably will be more efficient than in the discussion below. This issue deserves further attention.

Let  $s_m^2$  be the sample variance of  $m$  errors  $\{\epsilon_1, \dots, \epsilon_m\}$ . It is easily shown by calculations analagous to those of section 5.1 that replacing absolute residuals by sample standard deviations has the effect of changing the asymptotic covariance matrices (5.1), (5.4) and (5.5) to

$$\begin{aligned}
 \text{Pseudo-likelihood} &: \{(2 + \kappa) + 2/(m - 1)\} \{4N \xi(\beta, \theta)\}^{-1}; \\
 (5.7) \quad \text{Logarithm method} &: m \text{ var } \{ \log(s_m^2) \} \{4N \xi(\beta, \theta)\}^{-1}; \\
 \text{Weighted absolute residuals} &: \{m \delta_{\star} / (1 - \delta_{\star})\} \{N \xi(\beta, \theta)\}^{-1},
 \end{aligned}$$

where  $\delta_{\star} = \text{var}(s_m)$ . Table 3 contains the asymptotic relative efficiencies of using sample standard deviations to using transformations of absolute residuals for various values of  $m$  when the errors are standard normal. The values in the table for  $T(x) = x^2$  and  $x$  indicate that if the data are approximately normally distributed, using sample standard deviations can entail a loss in efficiency with respect to using residuals if  $m$  is small. For substantial replication ( $m \geq 10$ ), using sample standard deviations produces a slight edge in efficiency with respect to weighted absolute residuals for  $T(x) = x$ .

The second column of Table 3 shows that for the logarithm method, using sample standard deviations surpasses using residuals in terms of efficiency except when  $m = 2$  and is more than twice as efficient for large  $m$ . In its raw form,  $\log |r_i|$  is very unstable because, at least occasionally,  $|r_i| \approx 0$ .

producing a wild "outlier" in the regression. The effect of using sample standard deviations is to decrease the possibility of such inliers; the sample standard deviations will be likely more uniform, especially as  $m$  increases. The implication is that the logarithm method should not be based on residuals unless remedial measures are taken. The suggestion to trim a few of the smallest absolute residuals before using this method is clearly supported by the theory; presumably, such trimming would reduce or negate the theoretical superiority of using sample standard deviations.

Table 4 contains the asymptotic relative efficiencies of weighted squared sample standard deviations and logarithms of these to weighted squared residuals under normality of the errors. The first column is the efficiency of Raab's method to pseudo-likelihood, and the second column is the efficiency of the Rodbard and Frazier method to pseudo-likelihood. The results of the table imply that using the Raab and Rodbard and Frazier methods, which are popular in the analysis of radioimmunoassay data, can entail a loss of efficiency when compared to methods based on weighted squared residuals. Davidian (1986) has shown that the Rodbard and Frazier estimator can have a slight edge in efficiency over the weighted squared residuals methods for some highly contaminated normal distributions. From (5.7), the squared residual methods will be more efficient than Raab's method in the limit. Also note that the entries for  $T(x) = x$  and  $\log x$  in Table 3 for  $m = \infty$  are the reciprocals of the first row of Table 2 and that the entries for last row of Table 4 are 1.0; thus if both  $N$  and  $m$  grow large all the methods yield the same results.

Table 4 also addresses the open question as to whether Raab's method is asymptotically more efficient than the Rodbard and Frazier method for normally distributed data. The answer is a general yes, thus agreeing with the Monte-Carlo evidence available when the variance is a power of the mean. The results of this section suggest that in the case of assay data containing pairs

for which only  $Y_{ij}$  is observed, an estimator for  $\theta$  combining estimation based on residuals for the observations for which  $x_i$  is known and on standard deviations otherwise in an appropriately weighted fashion would offer some improvement over the methods currently employed; see Carroll, Davidian and Smith (1986).

## 6. DISCUSSION

In Section 4 we constructed a general theory of regression-type estimation for  $\theta$  in the heteroscedastic model (1.1). This theory includes as special cases common methods described in Section 3 and allows for the regression to be based on absolute residuals from the current regression fit as well as sample standard deviations in the event of replication at each design point. Under various restrictions such as symmetry or small  $\sigma$ , when the variance function  $g$  does not depend on  $\beta$ , we showed in Sections 4 and 5 that we can draw general conclusions about this class of estimators as well as make comparisons among the various methods.

When employing methods based on residuals, one should weight the residuals appropriately and iterate the process. There can be large relative differences among the methods in terms of efficiency. Under symmetry of the errors, squared residuals are preferable for approximately normally distributed data, but this preference is tenuous, these can be highly non-robust under only slight departures from normality; methods based on logarithms or the absolute residuals themselves exhibit relatively more robust behavior. For the small amount of replication found in practice, using sample standard deviations rather than residuals can entail a loss in efficiency if estimation is based on the squares of these quantities or the quantities themselves. For the

logarithm method based on residuals, trimming the smallest few absolute residuals is essential, since for normal data using sample standard deviations is almost always more efficient than using residuals, even for a small number of replicates. Popular methods applications such as radioimmunoassay based on sample means and sample standard deviations can be less efficient than methods based on weighted squared residuals. In some instances, the precision with which we can estimate  $\theta$  depends on the relative range of values of the mean responses, not their actual values, so that immediate implications for design are suggested.

Efficient variance function estimation in heteroscedastic regression analysis is an important problem in its own right. There are important differences in estimators for variance when it is modeled parametrically.

#### REFERENCES

- Abramowitz, M. and Stegun, I. A. (1972). Handbook of Mathematical Functions. Dover Publications, New York.
- Amemiya, T. (1977). A note on a heteroscedastic model. Journal of Econometrics 6, 365-370 and corrigenda 8, 265.
- Beal, S. L. and Sheiner, L. B. (1985). Heteroscedastic nonlinear regression with pharmacokinetic type data. Preprint.
- Box, G. E. P. (1986). Studies in quality improvement: signal to noise ratios, performance criteria and statistical analysis: part I. Center for Quality and Productivity Improvement, University of Wisconsin-Madison, Report #11.
- Box, G. E. P. and Hill, W. J. (1974). Correcting inhomogeneity of variance with power transformation weighting. Technometrics 16, 385-389.
- Box, G. E. P. and Meyer, R. D. (1986). Dispersion effects from fractional designs. Technometrics 28, 19-28.
- Box, G. E. P. and Ramirez, J. (1986). Studies in quality improvement: signal to noise ratios, performance criteria and statistical analysis: part II. Center for Quality and Productivity Improvement,

University of Wisconsin-Madison, Report #12.

- Carroll, R. J. (1982a). Adapting for heteroscedasticity in linear models. Annals of Statistics 10, 1224-1233.
- Carroll, R. J. (1987). The effects of variance function estimation on prediction and calibration: an example. Preprint.
- Carroll, R. J., Davidian, M. and Smith, W. (1986) Variance functions and the minimum detectable concentration in radioimmunoassay. Preprint.
- Carroll, R. J., and Ruppert, D. (1982b). A comparison between maximum likelihood and generalized least squares in a heteroscedastic linear model. Journal of the American Statistical Association 77, 878-882.
- Carroll, R. J., and Ruppert, D. (1982a). Robust estimation in heteroscedastic linear models. Annals of Statistics 10, 429-441.
- Carroll, R. J., and Ruppert, D. (1984). Power transformations when fitting theoretical models to data. Journal of the American Statistical Association 79, 321-328.
- Carroll, R. J., Wu, C. F. J. and Ruppert, D. (1987). Variance expansion and the bootstrap in generalized least squares. Preprint.
- Davidian, M. (1986). Variance function estimation in heteroscedastic regression models. Unpublished Ph.D. dissertation, University of North Carolina at Chapel Hill.
- Davidian, M. and Carroll, R.J. (1987). A note on extended quasi-likelihood. Preprint.
- Efron, B. (1986). Double exponential families and their use in generalized linear regression. Journal of the American Statistical Association 81, 709-721.
- Giltinan, D. M., Carroll, R. J. and Ruppert, D. (1986). Some new methods for weighted regression when there are possible outliers. Technometrics 28, 219-230.
- Glejser, H. (1969). A new test for heteroscedasticity. Journal of the American Statistical Association 64, 316-323.
- Gong, G. and Samaniego, F.J. (1981). Pseudo-maximum likelihood estimation: theory and applications. Annals of Statistics 9, 861-869.
- Harville, D. (1977). Maximum likelihood approaches to variance component estimation and to related problems. Journal of the American Statistical Association 79, 302-308.
- Harvey, A. C. (1976). Estimating regression models with multiplicative heteroscedasticity. Econometrica 44, 461-465.

- Huber, P. J. (1981). Robust Statistics. John Wiley and Sons, New York.
- Jacquez, J. A., Mather, F. J. and Crawford, C. R. (1968). Linear regression with non-constant, unknown error variances: sampling experiments with least squares and maximum likelihood estimators. Biometrics 24, 607-626.
- Jobson, J. D. and Fuller, W. A. (1980). Least squares estimation when the covariance matrix and parameter vector are functionally related. Journal of the American Statistical Association 75, 176-181.
- McCullagh, P. (1983). Quasi-likelihood functions. Annals of Statistics 11, 59-67.
- McCullagh, P. and Nelder, J. A. (1983). Generalized Linear Models. Chapman & Hall, New York.
- Nel, D. G. (1980). On matrix differentiation in statistics. South African Statistical Journal 14, 87-101.
- Nelder, J. A. and Pregibon, D. (1987). An extended quasi-likelihood function. Biometrika, to appear.
- Patterson, H. D. and Thompson, R. (1971). Recovery of inter-block information when block sizes are unequal. Biometrika 58, 545-554.
- Raab, G. M. (1981a). Estimation of a variance function, with application to radioimmunoassay. Applied Statistics 30, 32-40.
- Rodbard D. and Frazier, G. R. (1975). Statistical analysis of radioligand assay data. Methods of Enzymology 37, 3-22.
- Rosenblatt, J. R. and Spiegelman, C. H. (1981). Discussion of the paper by Hunter and Lamboy. Technometrics 23, 329-333.
- Rothenberg, T. J. (1984). Approximate normality of generalized least squares estimates. Econometrica 52, 811-825.
- Ruppert, D. and Carroll, R. J. (1980). Trimmed least squares estimation in the linear model. Journal of the American Statistical Association 77, 828-838.
- Sadler, W. A. and Smith, M. H. (1985). Estimation of the response-error relationship in immunoassay. Clinical Chemistry 31/11, 1802-1805.
- Stefanski, L. A. and Carroll, R. J. (1985). Covariate measurement error in logistic regression. Annals of Statistics 13, 1335-1351.
- Theil, H. (1971). Principles of Econometrics. New York: John Wiley and Sons.

Watters, R. L., Carroll, R. J. and Spiegelman, C. H. (1987). Error modeling and confidence interval estimation for inductively coupled plasma calibration curves. Preprint.

Williams, J. S. (1975). Lower bounds on convergence rates of weighted least squares to best linear unbiased estimators. In A Survey of Statistical Design and Linear Models, J. N. Srivastava, editor. Amsterdam, North Holland.

Wolter, K. M., and Fuller, W. A. (1982). Estimation of nonlinear errors-in-variables models. Annals of Statistics 10, 539-548.

#### APPENDIX A. PROOFS OF MAJOR RESULTS

We now present sketches of the proofs of the theorems of Section 4. Our exposition is brief and nonrigorous as our goal is to provide general insights. In what follows, we assume that

$$(A.1) \quad N^{1/2} \begin{bmatrix} \hat{\eta} - \eta \\ \hat{\theta} - \theta \end{bmatrix} = O_p(1);$$

under sufficient regularity conditions it is possible to prove (A.1). Such a proof would be long, detailed and essentially noninformative; see Carroll and Ruppert (1982a) for a proof of  $N^{1/2}$  consistency in a special case.

Sketch of proof of Theorem 4.1: From (4.1), a Taylor series, the fact that  $E [T\{d_1(\beta)\}] = M_1$  and laws of large numbers, we have

$$(A.2) \quad 0 = N^{-1/2} \sum_{i=1}^N (H_i/V_i) [T\{d_1(\hat{\beta}_*)\} - M_1(\hat{\eta}, \hat{\theta}, \hat{\beta}_*)] + o_p(1)$$

By the arguments of Ruppert and Carroll (1980) or Carroll and Ruppert (1982a),

$$\begin{aligned}
(A.3) \quad & N^{-1/2} \sum_{i=1}^N (H_i/V_i) [T\{d_i(\hat{\beta}_*)\} - T\{d_i(\beta)\}] \\
&= N^{-1/2} \sum_{i=1}^N (H_i/V_i) \dot{T}\{d_i(\beta)\} \{d_i(\hat{\beta}_*) - d_i(\beta)\} + o_p(1) \\
&= B_{3,N} N^{1/2} (\hat{\beta}_* - \beta) + o_p(1).
\end{aligned}$$

Applying this result to (A.2) along with a Taylor series in  $M_i$  gives

$$\begin{aligned}
0 &= N^{-1/2} \sum_{i=1}^N C_i + (B_{2,N} + B_{3,N}) N^{1/2} (\hat{\beta}_* - \beta) \\
&\quad - B_{1,N} N^{1/2} \begin{bmatrix} \hat{\eta} - \eta \\ \hat{\theta} - \theta \end{bmatrix} + o_p(1),
\end{aligned}$$

which is (4.2). □

Theorem 4.2 follows by a similar argument; in this case the representation (A.3) is unnecessary.

Sketch of proof of Theorem 4.3: We consider Theorem 4.2; the proof for Theorem 4.1 is similar. Recall here that (3.4) holds. In the following, all derivatives are with respect to the mean  $\mu_i$  and the definitions of  $C_i$  and  $M_i$  are as in (4.4).

Assumption (iv) implies that  $N^{1/2} \max_{1 \leq i \leq N} |\bar{Y}_i - \mu_i| \xrightarrow{P} 0$  so that a Taylor series in  $\eta$ ,  $\theta$  and  $\bar{Y}_i$  gives

$$(A.4) \quad B_{1,N} N^{1/2} \begin{bmatrix} \hat{\eta} - \eta \\ \hat{\theta} - \theta \end{bmatrix} = N^{-1/2} \sum_{i=1}^N C_i - N^{-1/2} \sum_{i=1}^N (\dot{M}_i H_i/V_i) (\bar{Y}_i - \mu_i)$$



$$+ N^{-1/2} \sum_{i=1}^N \{ (\dot{H}_1 / V_1) - (\dot{V}_1 / V_1) \} (\bar{Y}_{1.} - \mu_1) + o_p(1).$$

Since  $\bar{Y}_{1.} - \mu_1 = \sigma g(\mu_1, z_1, \theta) \bar{\epsilon}_{1.} \approx \lambda N^{-1/2} g(\mu_1, z_1, \theta) \bar{\epsilon}_{1.}$ , where  $\bar{\epsilon}_{1.}$  is the mean of the errors at  $x_1$ , we can write the last two terms on the right-hand side of (A.4) as

$$(A.5) \quad \lambda N^{-1} \sum_{i=1}^N \bar{\epsilon}_{1.} (q_{1,1} + q_{1,2} C_i)$$

for constants  $\{q_{1,j}\}$ . By assumption (v), since  $\bar{\epsilon}_{1.}$  has mean zero, (A.5) converges in probability to zero if  $E(\bar{\epsilon}_{1.} C_i) = 0$ , which holds under the assumption of symmetry. Thus, (A.5) converges to zero which from (A.4) completes the proof. Note that if we drop the assumption of symmetry, from (A.5) the asymptotic normal distribution of  $N^{1/2}(\hat{\theta} - \theta)$  will have mean

$$p\text{-}\lim_{N \rightarrow \infty} \{ \lambda B_{1,N}^{-1} N^{-1} \sum_{i=1}^N (\bar{\epsilon}_{1.} C_i q_{1,2}) \}.$$

□

## APPENDIX B. CHARACTERIZATION OF RESTRICTED MAXIMUM LIKELIHOOD

Let  $\hat{\beta}_*$  be a generalized least squares estimator for  $\beta$ . Assume first that  $g$  does not depend on  $\beta$ . Let the prior distribution for the parameters  $\pi(\beta, \theta, \sigma)$  be proportional to  $\sigma^{-1}$ . The marginal posterior for  $\theta$  is hard to compute in closed form for nonlinear regression. Following Box and Hill (1974) and Beal and Sheiner (1986), we have the linear approximation

$$f(x_1, \beta) \approx f(x_1, \hat{\beta}_*) + f_{\beta}(x_1, \hat{\beta}_*)^t (\beta - \hat{\beta}_*).$$

Replacing  $f(x_1, \beta)$  by its linear expansion, the marginal posterior for  $\theta$  is proportional to

$$(B.1) \quad p(\theta) = \frac{\{\prod_{i=1}^N g_i^2(\theta)\}^{-1/2}}{\hat{\sigma}_G^{(N-p)}(\theta) \{\text{Det } S_G(\theta)\}^{1/2}}, \text{ where}$$

$$\hat{\sigma}_G^2(\theta) = (N-p)^{-1} \sum_{i=1}^N r_i^2 / g^2(z_i, \hat{\beta}_*, \theta).$$

$$S_G(\theta) = N^{-1} \sum_{i=1}^N f_{\beta}(x_i, \hat{\beta}_*) f_{\beta}(x_i, \hat{\beta}_*)^t / g^2(z_i, \hat{\beta}_*, \theta).$$

and where  $\text{Det } A$  = determinant of  $A$ . If the variances depend on  $\beta$ , we extend the Bayesian arguments by replacing  $g_i(\theta)$  by  $g(z_i, \hat{\beta}_*, \theta)$ .

Let  $H$  be the hat matrix  $H$  evaluated at  $\hat{\beta}_*$  and let  $h_{ii} = h_{ii}(\hat{\beta}_*, \theta)$ . From (3.1), pseudo-likelihood solves in  $(\theta, \sigma)$

$$(B.2) \quad \sum_{i=1}^N [r_i^2 / \{\sigma^2 g^2(z_i, \hat{\beta}_*, \theta)\}] \begin{bmatrix} 1 \\ v_{\theta}(z_i, \hat{\beta}_*, \theta) \end{bmatrix} = \sum_{i=1}^N \begin{bmatrix} 1 \\ v_{\theta}(z_i, \hat{\beta}_*, \theta) \end{bmatrix}.$$

Since  $H$  is idempotent, the left hand side of (B.2) has approximate expectation

$$(B.3) \quad \sum_{i=1}^N \begin{bmatrix} 1 - p/N \\ v_{\theta}(z_i, \hat{\beta}_*, \theta) (1 - h_{ii}) \end{bmatrix}$$

To modify pseudo-likelihood to account for loss of degrees of freedom, equate the left hand side of (B.2) to (B.3). From matrix computations as in Nel (1980), this can be shown to be equivalent to restricted maximum likelihood.

Table 1

Description of Some Methods for Variance Function Estimation

<u>Maximum Likelihood</u>	Normal theory maximum likelihood in $\beta$ , $\sigma$ , $\theta$ .
<u>Pseudo-likelihood</u>	Normal theory maximum likelihood when $\beta$ is set to current value. When iterated, equivalent to maximum likelihood if the variance does not depend on $\beta$ .
<u>Weighted Squared Residuals</u>	Regress squared residuals on the variance, function, weight inversely with squared current variance estimate.
<u>Weighted Absolute Residuals</u>	Regress absolute residuals on the standard deviation function, weight inversely with current variance estimate.
<u>Logarithm Method</u>	Regress logarithm of absolute residuals on log of standard deviation function. Be wary of near zero residuals.
<u>Restricted Maximum Likelihood</u>	Pseudo-likelihood corrected for leverage. Maximizes marginal posterior for noninformative prior.
<p>All of the preceding except Restricted Maximum Likelihood have analogues formed by replacing absolute residuals by sample standard deviations in the case of replication. The following are based on the mean function or design being fully or partially unknown and are often used in assays.</p>	
<u>Rodbard and Frazier</u>	Regress log sample standard deviation on log sample mean, where the variance function depends on $\beta$ only through the means.
<u>Modified Maximum Likelihood</u>	Modified functional maximum likelihood (equation (2.5)), where variance function depends on $\beta$ only through means.
<u>Sadler and Smith</u>	Same as Modified Maximum Likelihood, but means estimated by sample means.

Table 2

Asymptotic relative efficiency of appropriately weighted regression methods based on a function  $T$  of absolute residuals and the method based on logarithms of absolute residuals with respect to appropriately weighted regression methods based on squared residuals for underlying contaminated normal error distributions with distribution function  $F(x) = (1 - \alpha)\phi(x) + \alpha\phi(x/3)$ .

$T(x)$

contamination fraction $\alpha$	<u><math>x</math></u>	<u><math>x^{2/3}</math></u>	<u><math>x^{1/2}</math></u>	<u><math>x^{1/3}</math></u>	<u><math>\log x</math></u>
0.000	0.876	0.772	0.693	0.606	0.405
0.001	0.948	0.841	0.756	0.662	0.440
0.002	1.016	0.906	0.816	0.715	0.480
0.010	1.439	1.334	1.216	1.075	0.720
0.050	2.035	2.100	1.996	1.823	1.220

Table 3

Asymptotic relative efficiency of regression methods based on a function  $T$  of sample standard deviations relative to using regression methods based on a function  $T$  of absolute residuals under normality for  $T(x)$  (weighted methods).

<u>m</u>	<u><math>x^2</math></u>	<u><math>T(x)</math></u>	
		<u>log x</u>	<u>x</u>
2	0.500	0.500	0.500
3	0.667	1.000	0.696
4	0.750	1.320	0.801
$\vdots$	$\vdots$	$\vdots$	$\vdots$
9	0.889	1.932	0.986
10	0.900	1.984	1.001
$\infty$	1.000	2.467	1.142

Table 4

Asymptotic relative efficiency of regression methods based on a function  $T$  of sample standard deviations relative to regression methods based on weighted squared residuals under normal errors.

$m$	$\frac{x^2}{T(x)}$	$\log x$
2	0.500	0.203
3	0.667	0.405
4	0.750	0.535
5	0.800	0.620
6	0.833	0.680
7	0.857	0.723
8	0.875	0.757
9	0.889	0.783
10	0.900	0.804
$\infty$	1.000	1.000

END  
DATE  
FILMED  
JAN  
1988